

ECE 538 - Advanced Computer Architecture  
Mid-Term Exam (Individual)

90

Date: 10/27/2021 Time: 2:00PM - 3:30PM

Test Guidelines:

- Open-book, Open-note, No Internet.
- Calculators are permitted, but no communication devices of any kind are allowed.
- Show your answers in the space provided for them. Write neatly and be well organized.
- The test is due exactly at the end of the class period. Each question is marked with its number of points, use your time wisely.

Suggestion:

- Whenever possible, show your work and your thought process. This will make it easier for us to give you partial credit.

Please sign your name below to indicate that you have read the above and have followed all the rules in the Academic Integrity according to the UNM student handbook.

Name: David Kirby

ID: 101652098

Score:

1. (20 points) There is a designed 5-stage, pipelined processor and synthesized it for a 45nm process technology node with a target clock rate of 1GHz. During power analysis, the processor is at the target clock rate with a supply voltage of 1.0V. This processor draws 70mW of dynamic power and 10mW of static power. Answer following questions by considering power and energy trade-offs:

- (a) (5 points) If a cryptographic operation takes 0.5 second to complete on the processor, what is the energy per operation at the target clock rate assuming back-to-back encryption operations?

$$\begin{aligned} \text{Power}_{\text{total}} &= \text{Power}_{\text{dynamic}} + \text{Power}_{\text{static}} \\ &= 70\text{mW} + 10\text{mW} \\ &= 80\text{mW} \end{aligned}$$

5

$$\begin{aligned} \text{Energy per operation} &= \text{Power}_{\text{total}} \times \text{time to complete operation} \\ &= 80\text{mW} \times 0.5\text{s} \\ &= 0.04\text{ J/op} \end{aligned}$$

- (b) (5 points) For common cryptographic operations on the processor, if you were to slow the clock down to 500MHz without adjusting the voltage, what would be the energy per operation? What would be the overall power draw?

$$\begin{aligned} \text{Power}_{\text{dynamic}}' &= \text{Power}_{\text{dynamic}} \times \frac{\text{new clk rate}}{\text{original clk rate}} \quad \text{Power}_{\text{static}} \text{ stay same!} \\ &= 70\text{mW} \times \frac{500\text{MHz}}{1\text{GHz}} \\ &= 35\text{mW} \end{aligned}$$

5

$$\begin{aligned} \text{Power}_{\text{total}}' &= \text{Power}_{\text{dynamic}}' + \text{Power}_{\text{static}} \\ &= 35\text{mW} + 10\text{mW} \\ &= 45\text{mW} \end{aligned}$$

$$\begin{aligned} \text{Energy per operation} &= \text{Power}_{\text{total}}' \times \text{time to complete operation} \\ &= 45\text{mW} \times 0.5\text{s} \\ &= 0.0225\text{ J/op} \end{aligned}$$

- (c) (5 points) If safely drop the voltage to 0.5V when operating at a 500MHz clock, recalculate the power draw and energy per operation. Assume the leakage current remains the same.

$$\begin{aligned} \text{Power}_{\text{dynamic}}' &= P_{\text{dynamic}} \times \frac{\text{new clk rate}}{\text{original clk rate}} \times \left( \frac{\text{new voltage}}{\text{old voltage}} \right)^2 \\ &= 70\text{mW} \times \left( \frac{500\text{MHz}}{1\text{GHz}} \right) \times \left( \frac{0.5\text{V}}{1.0\text{V}} \right)^2 \\ &= 8.75\text{mW} \end{aligned}$$

5

$$\text{Power}_{\text{total}}' = \text{Power}_{\text{dynamic}}' + \text{Power}_{\text{static}} = 13.75\text{mW}$$

$$\begin{aligned} \text{Power}_{\text{static}}' &= \text{Power}_{\text{static}} \times \frac{\text{new volt}}{\text{old volt}} \\ &= 10\text{mW} \times \left( \frac{0.5\text{V}}{1\text{V}} \right) \\ &= 5\text{mW} \end{aligned}$$

$$\text{Energy per operation} = P_{\text{total}}' \times \text{time to complete} = 13.75\text{mW} \times 0.5\text{s} = 0.006875\text{ J/op}$$

- (d) (5 points) Assuming your system performs one operation every second and gates the clock off in between when not performing a cryptographic operation, what would be the energy per operation? Also assume the original 1GHz clock rate and 1.0V supply voltage.

Power<sub>dynamic</sub>:

?

This is the dynamic voltage frequency scaling technique

2. (10 points) There is one cluster supported by the Center for Advanced Research Computing (CARC) at UNM. Assume the cluster has 500 computers, each of them with a MTTF of 25 days, and the failures follow an exponential distribution and are independent.

- (a) (5 credits) If  $1/5$  of the computers fail, the cluster is considered to fail. What is the MTTF of the cluster? Under this failure model, does adding more computers increase the MTTF on the cluster? Why?

$$\text{MTTF} = \frac{25 \text{ days}}{500 \text{ computers}} \times \frac{1}{5} \times 500 = 5 \text{ days}$$

5 Under this failure model, MTTF is not dependent on the number of computers and therefore adding more would not have an effect on the MTTF of the cluster since we are monitoring a  $\frac{1}{5}$  ratio

- (b) (5 credits) For the same amount of money, one could buy 800 computers, each with MTTF of 20 days. Assume that the cluster (implementation with either 800 less reliable computers or 500 original computers) is considered to fail if a single computer fails. Repairing the less reliable cluster configuration is 10% less expensive. Which cluster would be better?

$$\text{MTTF} = \frac{25}{500} = \frac{1}{20} \text{ days} = 1.2 \text{ hours}$$

$$\text{MTTF}' = \frac{20}{800} = \frac{1}{40} \text{ days} = 0.6 \text{ hours}$$

5 Failure rate is doubled and even with 10% less expensive equipment, I would recommend the older, more reliable cluster

3. (15 points) Define and briefly describe the types of cache misses (also called the three C's) discussed in lecture. List other type(s) of cache misses as you know.

1) Compulsory: these are initial misses due to an empty (cold) cache

15 2) Conflict: these are the misses due to a rigid block placement strategy (i.e., low associativity)

3) capacity: These are misses due to the cache being too small to hold the entire working set of data and instructions

others) coherency: misses due to the cache coherence protocol used for sharing memory amongst processors



4. (20 points) Suppose that in 1000 memory references there are 25 misses in the L1 cache and 12 misses in the L2 cache.

(a) (5 points) Calculate the local miss rate and the global miss rate for the L1 and L2 caches.

L1 Local miss rate =  $\frac{\# \text{ misses }}{\# \text{ accesses }} = \frac{25}{1000} = \frac{1}{40} = 0.025$  or 2.5%  
 Global miss rate = local miss rate for L1 =  $\frac{25}{1000} = \frac{1}{40} = 0.025$  or 2.5%

4 L2 Local miss rate =  $\frac{\# \text{ misses }}{\# \text{ accesses }} = \frac{12}{1000} = \frac{3}{250} = 0.012$  or 1.2%  
 Global miss rate = Local<sub>L1</sub> × Local<sub>L2</sub> =  $0.025 \times 0.012 = 0.0003$  or 0.03%

- (b) (5 points) Assume the miss penalty from L2 cache to main memory is 45 clock cycles, the hit time of L2 cache is 10 clock cycles, and the hit time of L1 is 1 clock cycle. Calculate the average memory access time.

AMAT = Hit Time<sub>L1</sub> + Miss rate<sub>L1</sub> × (hit time<sub>L2</sub> + miss rate<sub>L2</sub> × miss penalty<sub>L2</sub>)  
 $= 1cc + 0.025 \times (10cc + 0.012 \times 45cc)$   
 $= 1.025 \times (10.54cc)$   
 $= 10.8035cc$

- (c) (5 points) Assume this is an in order pipeline. If there are 1.25 memory references per instruction, what are the average stall cycles per instruction?

3  $\frac{\text{stalls}}{\text{instr.}} = \frac{\text{misses}_{L1}}{\text{instr.}} \times \text{hit time}_{L2} + \frac{\text{misses}_{L2}}{\text{instr.}} \times \text{miss penalty}_{L2}$   
 $= \frac{0.025}{1.25} \times 10cc + \frac{0.012}{1.25} \times 45cc$   
 $= 0.2 + 0.432 = 0.632 \frac{cc}{\text{instr.}}$

- (d) (5 points) Briefly describe the principles of the temporal and spacial locality. Explain how they relate to caching.

5 Spatial locality states that if a data location is referenced, there is a high likelihood that nearby locations will be referenced as well. Temporal deals with time and states that if a data location is referenced, there is a high likelihood that it will be referenced again. Instruction cache uses lots of spatial and temporal locality compared to data access, and helps to speed up the memory hierarchy.

5. (15 points) List the techniques to reduce the cache miss rate, miss penalty, and hit time with brief discussion (e.g., the effect of too large block size).

Reduce miss rate by:

- larger block sizes: spatial locality but large block sizes
- larger caches: but longer hit time Suffer from conflict misses
- higher associativity
- multi-level caches

Reduce miss penalty by:

- giving priority to read misses over writes

Reduce hit time by:

- critical word first, merging write buffers, and prefetching
- avoiding address translation during indexing and small, simple 1st level cache.

6. (20 points) Describe how write buffers are used to improve the performance of the memory hierarchy. Indicate where the write buffers exist in the memory hierarchy. Use diagram to help with the description.

Write buffers improve performance by reducing cache misses and reducing unnecessary stalling on writes, especially for write-through.

20

Write buffers can exist essentially anywhere in the memory hierarchy, for example between the L1 and L2 caches as shown below. Also between memory and L2.

